

CHAPTER V

NUMBER SYSTEMS AND ARITHMETIC

- Decimal number expansion

$$73625_{10} = (7 \times 10^4) + (3 \times 10^3) + (6 \times 10^2) + (2 \times 10^1) + (5 \times 10^0)$$

- Binary number representation

$$10110_2 = (1 \times 2^4) + (0 \times 2^3) + (1 \times 2^2) + (1 \times 2^1) + (0 \times 2^0) = 22_{10}$$

- Hexadecimal number representation

$$\begin{aligned} 3E4B8_{16} &= (3 \times 16^4) + (14 \times 16^3) + (4 \times 16^2) + (11 \times 16^1) + (8 \times 16^0) \\ &= 255160_{10} \end{aligned}$$

Radix-10 Representation

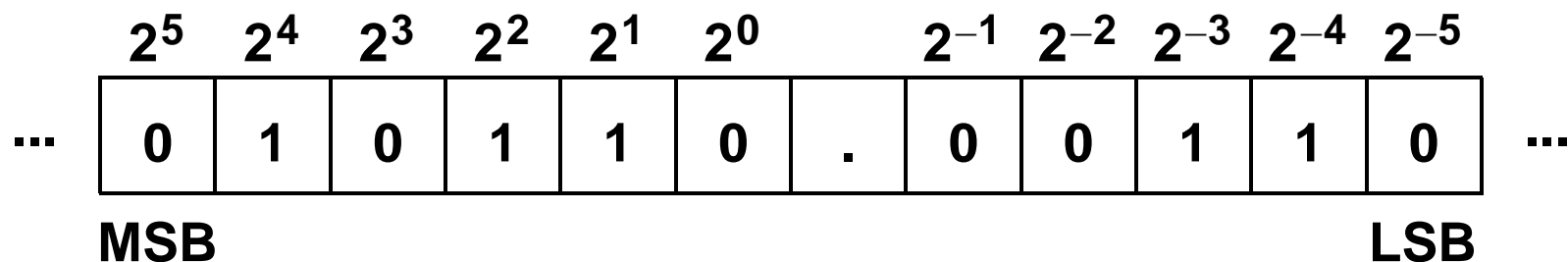
73625.4385₁₀

	10⁵	10⁴	10³	10²	10¹	10⁰		10⁻¹	10⁻²	10⁻³	10⁻⁴	10⁻⁵	
...	0	7	3	6	2	5	.	4	3	8	5	0	...

$$\begin{aligned} 73625.4385_{10} = & (7 \times 10^4) + (3 \times 10^3) + (6 \times 10^2) + (2 \times 10^1) + (5 \times 10^0) \\ & + (4 \times 10^{-1}) + (3 \times 10^{-2}) + (8 \times 10^{-3}) + (5 \times 10^{-4}) \end{aligned}$$

Radix-2 Representation

10110.0011₂



$$\begin{aligned} 10110.0011_2 &= (1 \times 2^4) + (0 \times 2^3) + (1 \times 2^2) + (1 \times 2^1) + (0 \times 2^0) \\ &\quad + (0 \times 2^{-1}) + (0 \times 2^{-2}) + (1 \times 2^{-3}) + (1 \times 2^{-4}) \\ &= 22.1875_{10} \end{aligned}$$

NUMBER SYSTEMS

OCTAL REPRESENTATION

Radix-8 Representation

26516.1731₈

	8^5	8^4	8^3	8^2	8^1	8^0		8^{-1}	8^{-2}	8^{-3}	8^{-4}	8^{-5}	
...	0	2	6	5	1	6	.	1	7	3	1	0	...

$$\begin{aligned} 26516.1731_8 &= (2 \times 8^4) + (6 \times 8^3) + (5 \times 8^2) + (1 \times 8^1) + (6 \times 8^0) \\ &\quad + (1 \times 8^{-1}) + (7 \times 8^{-2}) + (3 \times 8^{-3}) + (1 \times 8^{-4}) \\ &= 11598.24_{10} \end{aligned}$$

NUMBER SYSTEMS

HEXADECIMAL REPRES.

Radix-16 Representation

19AD6.F411₁₆

	16^5	16^4	16^3	16^2	16^1	16^0		16^{-1}	16^{-2}	16^{-3}	16^{-4}	16^{-5}	
...	0	1	9	A	D	6	.	F	4	1	1	0	...

$$\begin{aligned} 19AD6.F411_{16} &= (1 \times 16^4) + (9 \times 16^3) + (A \times 16^2) + (D \times 16^1) + (6 \times 16^0) \\ &\quad + (F \times 16^{-1}) + (4 \times 16^{-2}) + (1 \times 16^{-3}) + (1 \times 16^{-4}) \\ &\approx 105174.95_{10} \end{aligned}$$

NUMBER SYSTEMS

BINARY \leftrightarrow HEXADECIMAL

- NUMBER SYSTEMS
- BINARY REPRES.
- OCTAL REPRES.
- HEXADECIMAL REPRES.

BINARY \leftrightarrow HEXADECIMAL

$0000_2 = 0_{16}$	$1000_2 = 8_{16}$
$0001_2 = 1_{16}$	$1001_2 = 9_{16}$
$0010_2 = 2_{16}$	$1010_2 = 10 (A_{16})$
$0011_2 = 3_{16}$	$1011_2 = 11 (B_{16})$
$0100_2 = 4_{16}$	$1100_2 = 12 (C_{16})$
$0101_2 = 5_{16}$	$1101_2 = 13 (D_{16})$
$0110_2 = 6_{16}$	$1110_2 = 14 (E_{16})$
$0111_2 = 7_{16}$	$1111_2 = 15 (F_{16})$

BINARY \rightarrow HEXADECIMAL

Group binary by 4 bits from radix point

Examples:

$$\begin{array}{c} 0111 \ 1011_2 = 7B_{16} \\ \underbrace{\hspace{1cm}} \quad \underbrace{\hspace{1cm}} \\ 7 \quad \quad B \end{array}$$

$$\begin{array}{c} 10 \ 1010 \ 0110.1100 \ 01_2 = 2A6.C4_{16} \\ \underbrace{\hspace{1cm}} \quad \underbrace{\hspace{1cm}} \quad \underbrace{\hspace{1cm}} \quad \underbrace{\hspace{1cm}} \quad \underbrace{\hspace{1cm}} \\ 2 \quad \quad A \quad \quad 6 \quad \quad C \quad \quad 4 \end{array}$$

NUMBER SYSTEMS

BINARY \leftrightarrow OCTAL

- NUMBER SYSTEMS
- BINARY REPRES.
- OCTAL REPRES.
- BINARY \leftrightarrow HEXADECIMAL

BINARY \leftrightarrow OCTAL

$$000_2 = 0_8$$

$$001_2 = 1_8$$

$$010_2 = 2_8$$

$$011_2 = 3_8$$

$$100_2 = 4_8$$

$$101_2 = 5_8$$

$$110_2 = 6_8$$

$$111_2 = 7_8$$

BINARY \rightarrow OCTAL

Group binary bits by 3 from LSB

Examples:

$$\underbrace{10}_{2} \underbrace{100}_{4} \underbrace{110}_{6}_2 = 246_8$$

$$\underbrace{10}_2 \underbrace{101}_5 \underbrace{111}_7 \underbrace{011}_3 \underbrace{011}_3 \underbrace{11}_6_2 = 2573.36_8$$

NUMBER SYSTEMS

BINARY -> DECIMAL

- Perform radix-2 expansion
- Multiply each bit in the binary number by 2 to the power of its place.
Then sum all of the values to get the decimal value.

Examples:

$$10111_2 = (1 \times 2^4) + (0 \times 2^3) + (1 \times 2^2) + (1 \times 2^1) + (1 \times 2^0) = 23_{10}$$

$$\begin{aligned} 10110.0011_2 &= (1 \times 2^4) + (0 \times 2^3) + (1 \times 2^2) + (1 \times 2^1) + (0 \times 2^0) \\ &\quad + (0 \times 2^{-1}) + (0 \times 2^{-2}) + (1 \times 2^{-3}) + (1 \times 2^{-4}) \\ &= 22.1875_{10} \end{aligned}$$

NUMBER SYSTEMS

DECIMAL -> BINARY

- **Integer part:**
 - Modulo division of decimal integer by 2 to get each bit, starting with LSB.
- **Fraction part:**
 - Multiplication decimal fraction by 2 and collect resulting integers, starting with MSB.

Example: Convert 41.828125_{10}

$$\begin{array}{r}
 41 \bmod 2 = 1 \quad \text{LSB} \\
 20 \bmod 2 = 0 \\
 10 \bmod 2 = 0 \\
 5 \bmod 2 = 1 \\
 2 \bmod 2 = 0 \\
 1 \bmod 2 = 1 \quad \text{MSB}
 \end{array}$$

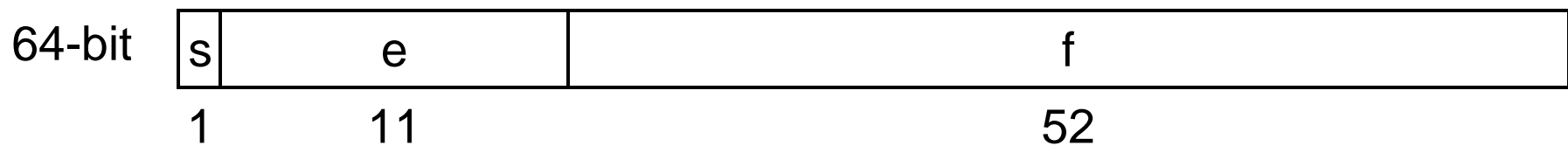
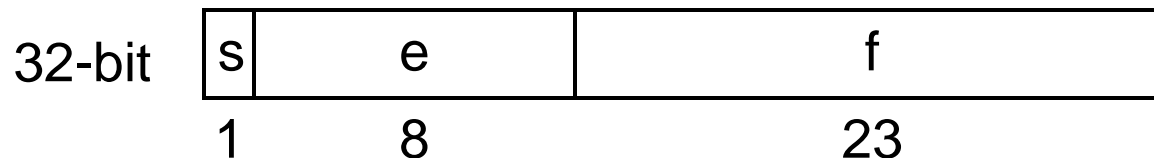
$$\begin{array}{r}
 0.828125 \times 2 = 1.65625 \quad \text{MSB} \\
 0.65625 \times 2 = 1.3125 \\
 0.3125 \times 2 = 0.625 \\
 0.625 \times 2 = 1.25 \\
 0.25 \times 2 = 0.5 \\
 0.5 \times 2 = 1.0 \quad \text{LSB}
 \end{array}$$

Therefore $41.828125_{10} = 101001.110101_2$

NUMBER SYSTEMS

FLOATING POINT NUMBERS

- Floating point numbers can be represented with a sign bit, a fraction (often referred to as the mantissa), and an exponent.
- Example 1: $-267.426 = -0.267426 \times 10^3$, where the sign is negative, the fraction is **0.267426** and the exponent is **3**.
- Example 2: $0101011.1001 = 0.1010111 \times 2^6$, where the sign is positive, the fraction is **0.1010111**, and the exponent is **0110**.
- Sample IEEE Floating-Point Formats



BINARY NUMBERS

UNSIGNED INTEGER

- The range for an n -bit radix- r unsigned integer is

$$[0, r_{10}^n - 1]$$

- Example: For a 16-bit binary unsigned integer, the range is

$$[0, 2^{16} - 1] = [0, 65535]$$

which has a binary representation of

$$\mathbf{0000\ 0000\ 0000\ 0000 = 0}$$

$$\mathbf{0000\ 0000\ 0000\ 0001 = 1}$$

$$\mathbf{0000\ 0000\ 0000\ 0010 = 2}$$

...

$$\mathbf{1111\ 1111\ 1111\ 1110 = 65534}$$

$$\mathbf{1111\ 1111\ 1111\ 1111 = 65535}$$

BINARY NUMBERS

SIGNED INTEGERS (1)

- The range for an n -bit radix- r signed integer is

$$[-r_{10}^{n-1}, r_{10}^{n-1} - 1]$$

- The most-significant bit is used as a sign bit, where **0** indicates a positive integer and **1** indicates a negative integer.

Example: For a 16-bit binary signed integer, the range is

$$[-2^{16-1}, 2^{16-1} - 1] = [-32768, 32767]$$



BINARY NUMBERS

SIGNED INTEGERS (2)

- There are a number of different representations for signed integers, each which has its own advantage
 - Signed-magnitude representation:
 - **1010 0001 0110 1111**
 - Signed-1's complement representation:
 - **1101 1110 1001 0000**
 - Signed-2's complement representation:
 - **1101 1110 1001 0001**
- The above examples are all the same number, **-8559_{10}** .

BINARY NUMBERS

SIGNED-MAGNITUDE

- The **signed-magnitude** binary integer representation is just like the **unsigned representation** with the addition of a **sign bit**.
- For instance, using 8-bits, the number -6_{10} can be represented as the 7-bit magnitude of 6_{10} using

000 0110

and then the sign bit appended to the MSB to form

1000 0110

BINARY NUMBERS

RADIX COMPLEMENTS

- The **radix complement**, or **r's complement**, of an integer representation for an n -digit integer is defined as

$$r^n_{10} - \text{number}_{10}$$

- The **diminished radix complement**, or **(r - 1)'s complement**, of an integer representation for an n -digit integer is defined as

$$(r^n_{10} - 1_{10}) - \text{number}_{10}$$

- Example: Find the r's and (r - 1)'s complement for **3764**₁₀

r's complement
 $10^5 - 3764 = 96236$

(r - 1)'s complement
 $(10^5 - 1) - 3764 = 96235$

BINARY NUMBERS

1'S COMPLEMENT

- The **1's complement** (diminished radix complement) binary integer representation for an n -bit integer is defined as

$$(2^n_{10} - 1_{10}) - \text{number}_{10}$$

- In essence, this takes the positive version of the number and flips all of the bits.
- For instance, using 8-bits, the number -6_{10} can be represented as the 8-bit positive number 6_{10} using

0000 0110

and then each of the bits flipped to form the 1's complement

1111 1001

BINARY NUMBERS

2'S COMPLEMENT

- The **2's complement** (radix complement) binary integer representation for an n -bit integer is defined as

$$2^n_{10} - \text{number}_{10}$$

- In essence, this takes the 1's complement and adds one.
 - For instance, using 8-bits, the number -6_{10} can be represented as the 8-bit positive number 6_{10} using

0000 0110

and then each of the bits flipped to form the 1's complement

1111 1001

and then add **1** to form the 2's complement

1111 1010

BINARY NUMBERS

SIGNED EXAMPLES

- Below are some examples for the signed binary numbers using 6 bits.

Decimal	Signed-magnitude	1's complement	2's complement
0	00 0000	00 0000	00 0000
1	00 0001	00 0001	00 0001
-1	10 0001	11 1110	11 1111
5	00 0101	00 0101	00 0101
-5	10 0101	11 1010	11 1011
12	00 1100	00 1100	00 1100
-12	10 1100	11 0011	11 0100
15	00 1111	00 1111	00 1111
-15	10 1111	11 0000	11 0001
16	01 0000	01 0000	01 0000
-16	11 0000	10 1111	11 0000

- Notice that all representations are the **same for positive numbers!!!!**

BINARY ARITHMETIC

UNSIGNED ADDITION

- Unsigned binary addition follows the standard rules of addition.
- Examples

$$\begin{array}{r}
 \mathbf{1111\ 0100} \text{ Carries} \\
 0011\ 1011 \\
 + 0111\ 1010 \\
 \hline
 1011\ 0101
 \end{array}$$

$$\begin{array}{r}
 \mathbf{0000\ 0010} \text{ Carries} \\
 1011\ 1001 \\
 + 0100\ 0101 \\
 \hline
 1111\ 1110
 \end{array}$$

$$\begin{array}{r}
 \mathbf{1111\ 0000} \text{ Carries} \\
 1111\ 1001 \\
 + 0100\ 1000 \\
 \hline
 \mathbf{1}\ 0100\ 0000
 \end{array}$$

$$\begin{array}{r}
 \mathbf{1110\ 0000\ 0000.0000} \text{ Carries} \\
 0101\ 1000\ 1001.1001 \\
 + 0011\ 0011\ 0100.01 \\
 \hline
 1000\ 1011\ 1101.1101
 \end{array}$$

BINARY ARITHMETIC

UNSIGNED SUBTRACTION

- Unsigned binary subtraction follows the standard rules.
- Examples

$$\begin{array}{r} \text{0000 0000} \text{ Borrows} \\ 1111 1001 \\ - 0100 1000 \\ \hline 1011 0001 \end{array}$$

$$\begin{array}{r} \text{1000 1000} \text{ Borrows} \\ 1011 1001 \\ - 0100 0101 \\ \hline 0111 0100 \end{array}$$

$$\begin{array}{r} \text{1000 0000} \text{ Borrows} \\ 0011 1011 \\ - 0111 1010 \\ \hline 1100 0001 \end{array}$$

$$\begin{array}{r} \text{0100 1110 1000.1000} \text{ Borrows} \\ 0101 1000 1001.1001 \\ - 0011 0011 0100.01 \\ \hline 0010 0101 0101.0101 \end{array}$$

BINARY ARITHMETIC

SIGNED ADDITION

- **Signed-magnitude**
 - Add magnitudes if signs are the same, give result the sign
 - Subtract magnitudes if signs are different. Absence or presence of an end borrow determines the resulting sign compared to the augend. If negative, then a 2's complement correction must be taken.
- **2's complement**
 - Add the numbers using normal addition rules. Carry out bit is discarded.
- **1's complement**
 - Easiest to convert to 2's complement, perform the addition, and then convert back to 1's complement. This is done as follows:
 - Add 1 to each integer, add the integers, subtract 1 from the result

BINARY ARITHMETIC

SIGNED SUBTRACTION

- Typically want to do addition or subtraction of **A** and **B** as follows.

$$\mathbf{SUM = A + B}$$

$$\mathbf{DIFFERENCE = A - B}$$

- If we use **2's complement**, we can make life easy on us since addition and subtraction are done in the same manner: **with addition only!!!**
- A subtraction can be re-represented as follows.

$$\mathbf{SUM = A + (-B)}$$

- Or in general any two numbers can be added as follows.

$$\mathbf{SUM = (\pm A) + (\pm B)}$$

BINARY ARITHMETIC

SIGNED MATH EXAMPLE

- Subtraction of signed numbers can best be done with 2's complement.
- Performed by taking the 2's complement of the subtrahend and then performing addition (including the sign bit).
 - Example:

$$\begin{array}{r}
 59 \\
 - 122 \\
 \hline
 \end{array}
 =
 \begin{array}{r}
 0011\ 1011 \\
 - 0111\ 1010 \\
 \hline
 \end{array}
 =
 \begin{array}{r}
 0011\ 1011 \\
 + 1000\ 0110 \\
 \hline
 1100\ 0001
 \end{array}
 = -(0011\ 1111) = -63$$

The diagram illustrates the subtraction of 122 from 59 using 2's complement. On the left, the decimal numbers are shown. An arrow labeled "2's complement" points from the subtrahend (122) to its 2's complement representation (0111 1010). The middle part shows the binary subtraction: 0011 1011 minus 0111 1010. The right part shows the addition of the 2's complement: 0011 1011 plus 1000 0110. The result is 1100 0001. An arrow labeled "discard carry out" points to the leading 1. Another arrow labeled "2's complement" points from the result 1100 0001 to its 2's complement, 0011 1111, which is then negated to get the final result -63. The word "Carries" is written in red next to the 0011 1100 part of the addition.